

MVAE Elektronikus Számolóközpont

Nagyhosszúságú aminosavszekvenciák közötti analógiák vizsgálata
számítógéppel

Farkas András

Az alábbiakban nagyhosszúságú aminosavszekvenciák közötti speciális analógia meghatározásával foglalkozunk.

A probléma a MTA SZBK Biokémiai Intézete Enzimológiai Részlegénél végzett fehérjeszerkezet-vizsgálatok kapcsán merült fel. Biokémiai megfogalmazása Dévényi Tibortól származik.

A probléma a következő: adott két vagy több aminosavszekvencia, keressük azt a leghosszabb elméleti aminosavszekvenciát, amelynek elemeit mindegyik adott aminosavszekvencia tartalmazza mégpedig, ugyanabban a sorrendben, mint amilyen sorrendben az elméleti aminosavszekvenciában szerepelnek.

Az aminosavszekvenciák közötti speciális analógiát az ilyen elméleti aminosavszekvencia létezése jelenti. A probléma matematikai megfogalmazása két adott szekvencia esetén a következő:

Adottak az $\{a_i\} = a_1, a_2, \dots, a_i, \dots, a_n$
 $\{b_i\} = b_1, b_2, \dots, b_i, \dots, b_m$ véges sorozatok.

A $\{c_k\} = c_1, c_2, \dots, c_k, \dots, c_l \quad l \leq \min(n, m)$

sorozatot $\{a_i\}$ és $\{b_i\}$ közös részsorozatának nevezzük, ha van olyan

$a_{i_1}, a_{i_2}, \dots, a_{i_k}, \dots, a_{i_l}$ részsorozata $\{a_i\}$ -nek,

továbbá van olyan

$b_{i_1}, b_{i_2}, \dots, b_{i_k}, \dots, b_{i_\ell}$ részsorozata $\{b_i\}$ -nek,

hogy $c_k = a_{i_k} = b_{i_k}$ minden $1 \leq k \leq \ell$ esetén.

Feladat: $\{a_i\}$ és $\{b_i\}$ sorozatokhoz megadni a leghosszabb közös részsorozatot, tehát amelyre ℓ maximális.

Két adott sorozathoz a leghosszabb közös részsorozat meghatározására szolgáló algoritmus a következő:

Egy $(n+1) \times (m+1)$ -es mátrix (melynek sorait és oszlopait 0-tól sorszámozzuk) $\lambda_{i,j}$ elemeit a következő módon értelmezzük:

1.) Legyen a 0. sor és 0. oszlop minden eleme 0:

$$\lambda_{0,j} = 0 \quad \lambda_{i,0} = 0 \\ 0 \leq j \leq m \quad 0 \leq i \leq n$$

2.) Értelmezzük a további elemeket pl. oszlopfolytonosan successzive:

$$\lambda_{i,j} = \begin{cases} 0 & \text{ha } a_i \neq b_j \\ \max_{\substack{0 \leq p < i \\ 0 \leq q < j}} (\lambda_{p,q}) + 1 & \text{ha } a_i = b_j \end{cases}$$

A maximális hosszúságú közös részsorozat, illetve részsorozatok a következő módon adódnak:

a maximális hosszúságú közös részsorozat hossza

$$\ell = \max_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} (\lambda_{i,j})$$

1.) Az utolsó eleme: $c_\ell = a_{i_\ell} = b_{j_\ell}$, ahol i_ℓ és j_ℓ az ℓ -el egyenlő $\lambda_{i,j}$ értékek közül valamelyiknek az indexei:

$$\lambda_{i_\ell, j_\ell} \in \left\{ \lambda_{i,j} : \lambda_{i,j} = \ell \right\}.$$

2.) Ha a k . eleme $c_k = a_{i_k} = b_{j_k}$, akkor a $k-1$. eleme

$$c_{k-1} = a_{i_{k-1}} = b_{j_{k-1}}$$

ahol i_{k-1} és j_{k-1} a $k-1$ értékkel egyenlő $\lambda_{i,j}$ elemek közül valamelyik olyannak az indexei, amelyre $i < i_k$ és $j < j_k$ fennáll.

$$\lambda_{i_{k-1}, j_{k-1}} \in \left\{ \lambda_{i,j} : \lambda_{i,j} = k-1, \begin{array}{l} i < i_k, \\ j < j_k \end{array} \right\}$$

Az algoritmus az összes lehetséges megoldást szolgáltatja.

A módszer általánosítható 2 dimenzióról n dimenzióra, ilyenkor n darab sorozat közös részsorozatának meghatározására alkalmas. A módszer helyességére vonatkozó konstruktív bizonyítás triviális.

Példa. (ld. 234. old.)

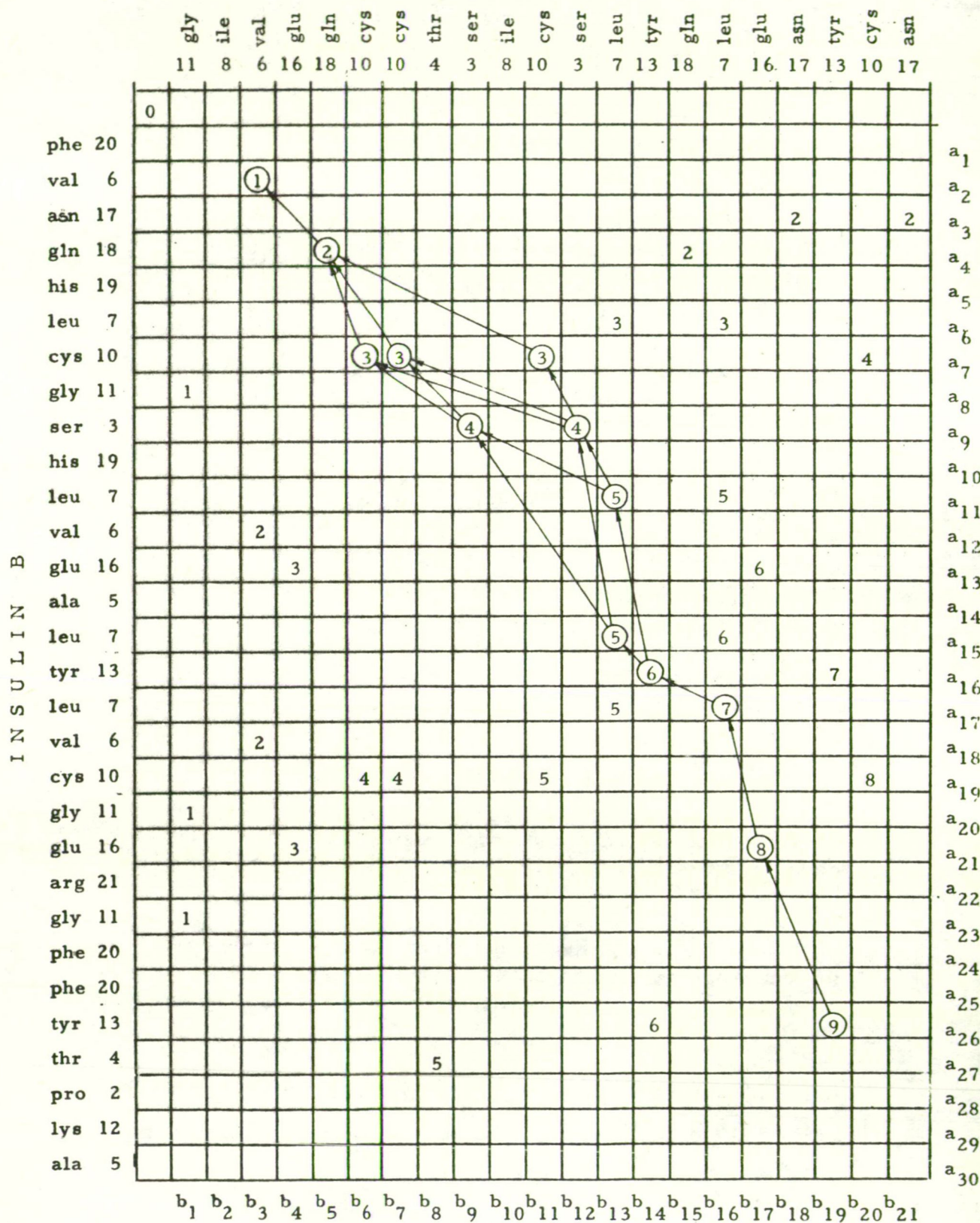
A módszert különböző nagyhosszúságú aminosavszekvenciák leghosszabb közös részszekvenciájának meghatározására alkalmaztuk.

Itt a szarvasmarha májból származó 506 tagú glutaminsav dehidrogenáz, illetve a rák izomból származó 333 tagú gliceraldehid 3 foszfát dehidrogenáz aminosavszekvenciák esetén kapott eredményeinket ismertetjük.

A leghosszabb közös részszekvencia 151 tagunak adódott.

A következőkben ismertetjük a két aminosavszekvenciát és

INSULIN A



1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
ALA	ASP	ARG	GLU	ASP	ASP	PRO	ASN	PHE	PHE	<u>LYS</u>	MET	VAL	GLU	<u>GLY</u>	PHE	<u>PHE</u>	ASP	ARG	<u>GLY</u>	ALA	SEN	<u>ILE</u>	VAL	GLU
26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
ASP	LYS	<u>LEU</u>	<u>VAL</u>	GLU	ASP	<u>LEU</u>	LYS	THR	ARG	GLN	THR	GLN	GLU	GLN	LYS	ARG	ASN	ARG	<u>VAL</u>	ARG	<u>GLY</u>	<u>ILE</u>	LEU	ARG
51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
ALA	GLN	HIS	SEN	HIS	GLN	ARG	THR	<u>PRO</u>	CYS	LYS	GLY	GLY	<u>ILE</u>	ARG	<u>TYR</u>	SEN	THR	ASP	<u>VAL</u>	<u>SEN</u>	<u>VAL</u>	ASP	<u>GLU</u>	<u>VAL</u>
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
<u>LYS</u>	ALA	LEU	ALA	SEN	<u>LEU</u>	<u>MET</u>	THR	TYR	LYS	CYS	ALA	<u>VAL</u>	<u>VAL</u>	ASP	VAL	PRO	PHE	GLY	<u>GLY</u>	ALA	<u>LYS</u>	ALA	GLY	<u>VAL</u>
101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125
<u>LYS</u>	<u>ILE</u>	ASN	PRO	LYS	ASN	TYR	THR	ASP	GLU	ASP	LEU	<u>GLU</u>	LYS	<u>ILE</u>	THR	ARG	THR	ARG	<u>PHE</u>	MET	GLU	LEU	THR	THR
126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150
ALA	MET	<u>GLU</u>	LEU	ALA	LYS	<u>LYS</u>	GLY	<u>PHE</u>	<u>ILE</u>	<u>GLY</u>	PRO	<u>GLY</u>	LEU	ASP	<u>VAL</u>	<u>PRO</u>	ALA	<u>PRO</u>	ASN	<u>MET</u>	SEN	THR	<u>GLY</u>	GLU
151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
ARG	<u>GLU</u>	MET	<u>SEN</u>	TRP	<u>ILE</u>	ALA	<u>ASP</u>	THR	TYR	ALA	<u>SEN</u>	THR	<u>ILE</u>	GLY	HIS	TYR	ASP	<u>ILE</u>	ASN	ALA	HIS	ALA	CYS	<u>VAL</u>
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200
THR	<u>LYS</u>	<u>PRO</u>	GLY	<u>ILE</u>	SEN	GLN	GLY	GLY	ILE	SEN	ALA	THR	GLY	ARG	<u>VAL</u>	<u>PHE</u>	GLY	ARG	<u>GLY</u>	<u>VAL</u>	<u>PHE</u>	GLY	HIS	<u>ILE</u>
201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225
<u>GLU</u>	ASN	<u>PHE</u>	<u>ILE</u>	GLU	ASN	ALA	SEN	TYR	MET	SEN	<u>ILE</u>	LEU	<u>GLY</u>	MET	THR	PRO	GLY	<u>PHE</u>	GLY	ASP	LYS	THR	<u>PHE</u>	ALA
226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250
<u>VAL</u>	<u>GLN</u>	GLY	<u>PHE</u>	GLY	ASN	<u>VAL</u>	<u>GLY</u>	LEU	HIS	<u>SEN</u>	MET	ARG	TYR	LEU	HIS	ARG	<u>PHE</u>	<u>GLY</u>	<u>ALA</u>	LYS	CYS	<u>VAL</u>	<u>ALA</u>	<u>VAL</u>
251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275
GLY	GLU	SEN	ASP	GLY	SEN	<u>ILE</u>	TRP	ASN	PRO	ASP	PRO	<u>ILE</u>	ASP	<u>PRO</u>	LYS	GLU	LEU	GLU	ASP	<u>PHE</u>	LYS	LEU	GLN	HIS
276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300
GLY	THR	<u>ILE</u>	LEU	<u>GLY</u>	<u>PHE</u>	<u>PRO</u>	LYS	ALA	<u>LYS</u>	<u>ILE</u>	TYR	GLU	GLY	SEN	<u>ILE</u>	LEU	<u>GLU</u>	<u>VAL</u>	ASP	<u>GLY</u>	ASP	<u>ILE</u>	LEU	<u>ILE</u>
301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325
PRO	ALA	ALA	SEN	GLU	<u>LYS</u>	GLN	<u>LEU</u>	THR	LYS	SEN	ASN	ALA	PRO	ARG	<u>VAL</u>	LYS	ALA	LYS	<u>ILE</u>	<u>ILE</u>	ALA	GLU	GLY	ALA
326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350
ASN	GLY	<u>PRO</u>	THR	<u>THR</u>	<u>PRO</u>	GLY	ALA	<u>ASP</u>	LYS	<u>ILE</u>	<u>PHE</u>	LEU	GLU	ARG	<u>ILE</u>	<u>ILE</u>	LYS	PRO	CYS	ASN	HIS	<u>VAL</u>	LEU	SEN
351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375
LEU	SEN	<u>PHE</u>	<u>PRO</u>	<u>ILE</u>	ARG	ARG	ASP	ASP	GLY	<u>SEN</u>	TRP	GLU	<u>VAL</u>	<u>ILE</u>	GLU	GLY	TYR	ARG	<u>ILE</u>	GLY	MET	<u>VAL</u>	<u>ILE</u>	<u>PRO</u>
376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400
<u>ASP</u>	LEU	TYR	<u>LEU</u>	ASN	ALA	GLY	GLY	<u>VAL</u>	<u>THR</u>	<u>VAL</u>	SEN	TYR	<u>PHE</u>	GLY	LEU	LYS	ASN	LEU	ASN	HIS	<u>VAL</u>	SEN	TYR	GLY
401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425
ARG	<u>LEU</u>	THR	<u>PHE</u>	<u>LYS</u>	TYR	<u>GLU</u>	ARG	ASP	<u>SEN</u>	ASN	TYR	HIS	LEU	LEU	<u>MET</u>	SEN	VAL	GLN	GLU	<u>SEN</u>	LEU	GLU	ARG	LYS
426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450
<u>PHE</u>	<u>GLY</u>	LYS	HIS	<u>GLY</u>	<u>GLY</u>	THR	<u>ILE</u>	PRO	<u>ILE</u>	<u>VAL</u>	PRO	<u>THR</u>	ALA	<u>GLU</u>	<u>PHE</u>	GLN	ASP	ARG	<u>ILE</u>	<u>SEN</u>	GLY	ALA	<u>SEN</u>	GLU
451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475
LYS	<u>ASP</u>	<u>ILE</u>	<u>VAL</u>	HIS	SEN	<u>GLY</u>	LEU	ALA	TYR	THR	MET	GLU	ARG	<u>SEN</u>	ALA	ARG	GLN	<u>ILE</u>	MET	ARG	THR	ALA	MET	<u>LYS</u>
476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500
TYR	ASN	LEU	<u>GLY</u>	LEU	ASP	<u>LEU</u>	ARG	THR	ALA	ALA	TYR	<u>VAL</u>	ASN	ALA	<u>ILE</u>	GLU	<u>LYS</u>	<u>VAL</u>	<u>PHE</u>	ARG	<u>VAL</u>	TYR	ASN	GLU
501	502	503	504	505	506																			
ALA	<u>GLY</u>	<u>VAL</u>	THR	<u>PHE</u>	THR																			

Glutaminsav dehidrogenáz.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
SER	LYS	ILE	GLY	ILE	ASP	GLY	PHE	GLY	ARG	ILE	GLY	ARG	LEU	VAL	LEU	ARG	ALA	ALA	LEU	SER	CYS	GLY	ALA	GLN
26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
VAL	VAL	ALA	VAL	ASN	ASP	PRO	PHE	ILE	ALA	LEU	GLU	TYR	MET	VAL	TYR	MET	PHE	LYS	TYR	ASP	SER	THR	HIS	GLY
51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
VAL	PHE	LYS	GLY	GLU	VAL	LYS	MET	GLU	ASP	GLY	ALA	LEU	VAL	VAL	ASP	GLY	LYS	LYS	ILE	THR	VAL	PHE	ASN	GLU
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
MET	LYS	PRO	GLU	ASN	ILE	PRO	TRP	SER	LYS	ALA	GLY	ALA	GLU	TYR	ILE	VAL	GLU	SER	THR	GLY	VAL	PHE	THR	THR
101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125
ILE	GLU	LYS	ALA	SER	ALA	HIS	PHE	LYS	GLY	GLY	ALA	LYS	LYS	VAL	VAL	ILE	SER	ALA	PRO	SER	ALA	ASP	ALA	PRO
126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150
MET	PHE	VAL	CYS	GLY	VAL	ASN	LEU	GLU	LYS	TYR	SER	LYS	ASP	MET	THR	VAL	VAL	SER	ASN	ALA	SER	CYS	THR	THR
151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
ASN	CYS	LEU	ALA	PRO	VAL	ALA	LYS	VAL	LEU	HIS	GLU	ASN	PHE	GLU	ILE	VAL	GLU	GLY	LEU	MET	THR	THR	VAL	HIS
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200
ALA	VAL	THR	ALA	THR	GLN	LYS	THR	VAL	ASP	GLY	PRO	SER	ALA	LYS	ASP	THR	ARG	GLY	GLY	ARG	GLY	ALA	ALA	GLN
201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225
ASN	ILE	ILE	PRO	SER	SER	THR	GLY	ALA	ALA	LYS	ALA	VAL	GLY	LYS	VAL	ILE	PRO	GLU	LEU	ASP	GLY	LYS	LEU	THR
226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250
GLY	MET	ALA	PHE	ARG	VAL	PRO	THR	PRO	ASP	VAL	SER	VAL	VAL	ASP	LEU	THR	VAL	ARG	LEU	GLY	LYS	GLU	CYS	SER
251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275
TYR	ASP	ASP	ILE	LYS	ALA	ALA	MET	LYS	THR	ALA	SER	GLU	GLY	PRO	LEU	GLN	GLY	PHE	LEU	GLY	TYR	THR	GLU	ASP
276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300
ASP	VAL	VAL	SER	SER	ASP	PHE	ILE	GLY	ASP	ASN	ARG	SER	SER	ILE	PHE	ASP	ALA	LYS	ALA	GLY	ILE	GLN	LEU	SER
301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325
LYS	THR	PHE	VAL	LYS	VAL	VAL	SER	TRP	TYR	ASP	ASN	GLU	PHE	GLY	TYR	SER	GLN	ARG	VAL	ILE	ASP	LEU	LEU	LYS
326	327	328	329	330	331	332	333																	
HIS	MET	GLN	LYS	VAL	ASP	SER	ALA																	

Gliceraldehyd 3 foszfát dehidrogenáz.

a megfelelő aminosavak aláhúzásával jelöljük azokat az aminosavakat, amelyek a leghosszabb közös részszekvenciába tartoznak.

A lehetséges megoldások közül egy "bal-alsó" megoldást választottunk.

Az eredményként adódó elméleti aminosavszekvenciát meg kell vizsgálnunk abból a szempontból, hogy hossza "kirívóan" nagy-e, azaz valószínűtlen esemény-e az, hogy a megvizsgált aminosavszekvenciákkal azonos hosszúságú és aminosavösszetételű, de véletlen aminosav-sorrendű szekvenciákban legalább az eredményben szereplő leghosszabb közös aminosav-szekvenciahossz adódik.

Ha valóban valószínűtlen esemény, akkor szükséges a "kirívó" hossz jelenségének biokémiai interpretációja.

A valószínűséget sztochasztikus szimulációval határozhatjuk meg, tekintettel arra, hogy a valószínűség explicit alakban való megadása problematikus.

